

CLAIMS

1. A method of generating a definition for a collection of source documents comprising:

identifying patterns common to each source document in the collection of source documents; and

constructing for an element type in the collection of source documents a restrictive general rule based on the identified common patterns.

2. The method of claim 1, wherein identifying common patterns comprises:

identifying common attribute names and types.

3. The method of claim 2, wherein identifying common patterns further comprises:

identifying restricted attribute values associated with the common attribute names and types.

4. The method of claim 2, wherein identifying common attribute names and types comprises:

determining the number of occurrences of each attribute name on an element type;

examining the attribute values for each occurrence of each attribute name on the same element type to determine the attribute type; and

determining if the attribute name occurs in association with the same attribute value on more than one element type.

1 5. The method of claim 3, wherein identifying restricted
2 attribute values comprises:

3 examining attribute values for each occurrence of an
4 attribute type in all of the source documents in the collection
5 of source documents; and

6 establishing an enumeration or a restricted range
7 appropriate to the attribute type.

1 6. The method of claim 5, wherein identifying restricted
2 attribute values further comprises:

3 applying a heuristic to identify errors in the
4 collection of source documents; and

5 adjusting the established enumeration or restricted
6 range for attribute values.

1 7. The method of claim 1, wherein constructing a
2 restricted general rule comprises:

3 constructing a content model that specifies the
4 sequence order and number of occurrences of sub-elements within
5 the common pattern.

1 8. The method of claim 2, wherein constructing a
2 restricted general rule comprises:

3 constructing attribute definitions and value rules for
4 each identified common attribute name and type.

1 9. The method of claim 1, further comprising:
2 identifying those patterns found to achieve a
3 predetermined threshold of commonness; and
4 constructing a restrictive general rule for those
5 identified patterns.

10. A computer program residing on a computer-readable medium for building a document type definition for a collection of source documents, the computer program comprising instructions causing a computer system to:

identify patterns common to each source document in the collection of source documents; and

construct for an element type in the collection of source documents a restrictive general rule based on the identified common patterns.

11. The computer program of claim 10, wherein the instructions to identify common patterns comprise instructions to:

identify common attribute names and types.

12. The computer program of claim 11, wherein the instructions to identify common patterns further comprise instructions to:

identify restricted attribute values associated with the common attribute names and types.

13. A computer system comprising:
a storage device for storing a set of source documents;
and

a computer processor configured by a document type definition building program to identify patterns common to each source document in the set of source documents and construct for an element type in the set of source documents a restrictive general rule base on the identified common patterns.

14. A method of converting a format of a first source document to a format of a similarly structured second source document, the method comprising:

identifying patterns common to the first and second source documents; and

using the identified common patterns to map elements and sub-elements in the first source document to equivalent elements and sub-elements in the second source document.

15. The method of claim 14, further comprising:

replacing tag names for each of the elements and sub-elements in the first source document with equivalent tag names of the elements and sub-elements in the second source document.

16. The method of claim 14, wherein identifying patterns common to the first and second source documents comprises:

examining document type definitions for the first and second source documents.

1 17. The method of claim 16, further comprising:
2 producing the document type definition for the first
3 source document if the document type definition for the first
4 source document does not already exist.

1 18. The method of claim 14, wherein identifying patterns
2 common to the first and second source documents comprises:
3 performing pattern matching.

1 19. The method of claim 14, wherein identifying patterns
2 common to the first and second source documents comprises:
3 matching heuristics of the patterns in the first source
4 document to heuristics of the patterns in the second source
5 document.

1 20. The method of claim 18, wherein identifying patterns
2 common to the first and second source documents further
3 comprises:
4 matching heuristics of the patterns in the first source
5 document to heuristics of the patterns in the second source
6 document.

7 21. The method of claim 14, wherein using uses the
8 identified common patterns to map automatically elements and sub-
9 elements in the first source document to equivalent elements and
10 sub-elements in the second source document.

1 22. A method of converting the format of a source document

2 to the format of a set of source documents, the set of source
3 documents having a structure similar to the first source
4 document, the method comprising:

5 identifying patterns common to the source document and
6 the set of source documents;

7 mapping elements and sub-elements in the common pattern
8 of the source document to equivalent elements and sub-elements
9 the common pattern of the set of source documents; and

10 replacing tag names for the each of the elements and
11 sub-elements in common pattern of the source document with the
12 equivalent tag names of the elements and sub-elements in common
13 pattern of the set of source documents.

23. The method of claim 22, wherein identifying patterns
common to the source document and the set of source documents
comprises:

4 examining document type definitions for the source
document and and the set of source documents.

24. The method of claim 23, further comprising:

2 producing the document type definition for the source
3 document if the document type definition for the source document
4 does not already exist.

1 25. A computer program residing on a computer-readable
2 medium for converting a format of a first source document to a
3 format of a similarly structured second source document, the
4 computer program comprising instructions causing a computer

5 system to:

6 identify patterns common to the first and second source
7 documents; and

8 use the identified common patterns to map elements and
9 sub-elements of the first source document to equivalent elements
10 and sub-elements of the second source document.

11 26. The computer program of claim 25, further comprising
12 instructions to:

13 replace tag names for the each of the elements and sub-
14 elements in the common pattern of the first source document with
15 equivalent tag names of the elements and sub-elements in the
16 common pattern of the second source document.

17 27. The computer program of claim 26, wherein the
18 instructions to identify patterns common to the source document
19 and the set of source documents comprise instructions to:

20 examine document type definitions for the source
21 document and and the set of source documents.

22 28. A computer system comprising:

23 a storage device for storing a source document and a
24 set of source documents, the source document having a format
25 different from that of the set of source documents; and

26 a computer processor configured by a mapping program to
27 identify patterns common to the source document and the set of
28 source documents and map elements and sub-elements in the common
29 pattern of the source document to equivalent elements and sub-

9 elements the common pattern of the set of source documents.

Add
A₂
Add
B₁